

Visual Profiling of Large Statistical Datasets

Martijn Tennekes
Edwin de Jonge, Piet Daas

January 31, 2011



Funded under Socio-economic Sciences & Humanities



Centraal Bureau voor de Statistiek

Outline

- Introduction
- Tableplot description
- Applications
- Implementation in R

Introduction

Large statistical dataset

- Administrative sources
- Survey data

Quality assessment at a technical level

Step 1: Technical checks (e.g. readability and convertability)

Step 2: Data profiling

- Representation and distribution of values
- Strange data patterns
- Occurrence of missing values

Introduction

Large statistical dataset

- Administrative sources
- Survey data

Quality assessment at a technical level

Step 1: Technical checks (e.g. readability and convertability)

Step 2: Data profiling

- Representation and distribution of values
- Strange data patterns
- Occurrence of missing values

Introduction

Large statistical dataset

- Administrative sources
- Survey data

Quality assessment at a technical level

Step 1: Technical checks (e.g. readability and convertability)

Step 2: Data profiling

- Representation and distribution of values
- Strange data patterns
- Occurrence of missing values

Introduction

Large statistical dataset

- Administrative sources
- Survey data

Quality assessment at a technical level

Step 1: Technical checks (e.g. readability and convertability)

Step 2: Data profiling

- Representation and distribution of values
- Strange data patterns
- Occurrence of missing values

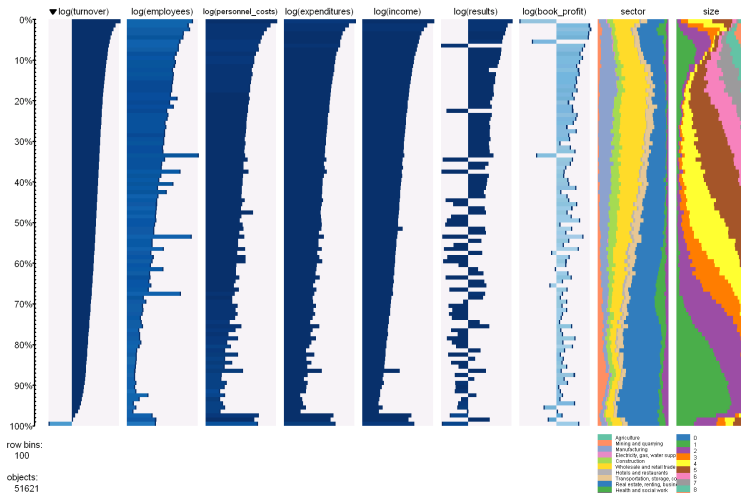
Introduction

Traditional approach:

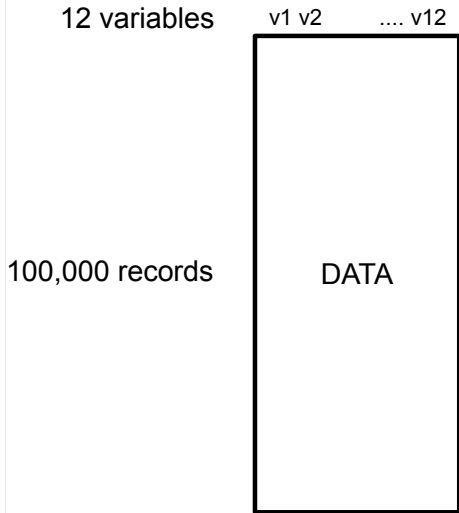
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	AC	AD
1	1	core	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
2	1	8.20 Ideal	J	IS2	61.5	55	326	335	336	243	645	55	2762	537	553	354	655	55	2875	63	925	411	636	57	2790	571	581	565	358	65
3	2	8.21 Premium	Z	IS1	59.8	61	326	339	334	231	653	55	2762	563	638	366	663	61	2875	608	833	401	612	59	2790	578	538	593	362	62
4	3	8.20 Good	J	VS1	59.9	60	329	336	334	230	653	60	2762	569	632	366	663	58	2875	614	836	404	614	58	2790	574	578	576	369	61
5	4	8.20 Premium	J	VS2	62.4	58	334	42	43	263	614	55	2763	603	597	371	626	60	2876	580	75	358	615	55	2790	581	586	596	368	62
6	5	8.21 Good	J	IS2	63.2	57	335	434	435	275	664	57	2763	589	638	370	666	53	2876	604	846	364	614	59	2790	576	576	368	61	
7	6	8.24 Very Good	J	VS2	62.8	57	336	394	386	248	644	57	2763	611	608	359	624	54	2877	574	576	357	618	60	2791	578	638	596	364	59
8	7	8.24 Very Good	J	VS1	62.5	57	336	395	338	247	657	60	2763	613	638	358	624	56	2877	575	57	357	627	58	2791	583	588	588	364	60
9	8	8.24 Very Good	J	VS1	61.9	56	337	407	411	253	62	56	2763	589	636	364	623	55	2877	58	833	362	615	58	2791	582	586	586	363	60
10	9	8.22 Fair	J	VS2	65.1	61	337	397	378	249	634	56	2764	584	588	3	627	56	2877	58	836	355	606	63	2792	583	578	571	371	62
11	10	8.23 Very Good	J	VS1	60.4	61	339	4	405	236	619	59	2764	589	572	353	601	62	2877	592	846	344	575	64	2792	584	537	368	63	
12	11	8.3 Good	J	IS1	64	55	339	425	428	273	628	60	2765	582	565	353	628	59	2878	586	838	369	613	58	2792	58	594	583	374	61
13	12	8.20 Good	J	VS1	62.8	58	340	333	333	246	613	60	2765	588	538	341	627	60	2878	582	839	367	627	62	2792	581	577	571	371	62
14	13	8.20 Premium	J	IS1	60.4	61	342	398	384	230	668	60	2765	582	565	357	607	58	2878	6	837	36	628	60	2792	582	584	584	369	62
15	14	8.31 Ideal	J	IS2	62.2	64	344	435	437	271	611	62	2765	574	577	368	615	55	2878	576	738	355	631	59	2792	578	573	563	61	61
16	15	8.20 Premium	J	IS2	62.2	62	345	378	375	227	618	59	2765	589	573	353	628	58	2878	588	538	377	608	57	2792	578	575	575	375	62
17	16	8.30 Premium	Z	IS1	60.9	58	345	438	442	260	612	57	2765	588	538	361	614	56	2878	581	838	358	627	59	2792	573	569	558	374	61
18	17	8.3 Ideal	J	IS2	62	54	348	431	434	286	618	56	2766	533	536	343	603	62	2879	556	535	354	628	60	2792	588	586	586	368	6
19	18	8.3 Good	J	IS1	61.4	64	351	423	428	277	668	57	2766	576	575	351	604	60	2879	572	57	345	611	61	2792	539	62	347	62	
20	19	8.3 Good	J	IS1	63.8	58	351	423	428	271	668	56	2766	589	581	33	627	58	2879	571	837	357	614	59	2792	588	585	585	374	61
21	20	8.3 Very Good	J	IS1	62.7	59	351	421	427	288	618	55	2767	589	572	352	601	57	2879	603	61	373	622	55	2792	574	570	570	371	62
22	21	8.3 Good	J	IS2	65.3	58	351	426	43	271	60	57	2767	58	567	35	623	55	2879	606	603	377	63	54	2794	589	573	573	374	61
23	22	8.23 Very Good	J	VS2	63.8	55	352	398	383	248	618	55	2767	574	576	354	603	54	2879	653	635	360	637	56	2794	585	582	575	375	61
24	23	8.23 Very Good	J	VS1	61.1	62	353	394	384	241	621	64	2767	582	585	33	608	60	2879	530	777	329	636	58	2797	588	586	586	361	61
25	24	8.31 Very Good	J	IS1	60.4	62	353	439	443	262	633	60	2768	562	581	352	602	56	2879	576	773	356	627	57	2797	571	575	575	378	61
26	25	8.31 Very Good	J	IS1	60.1	62	353	444	447	258	602	56	2768	583	587	352	607	61	2880	604	642	372	61	59	2797	584	585	584	371	62
27	26	8.23 Very Good	J	VS2	60.4	58	354	397	401	241	611	58	2768	586	573	348	639	57	2880	572	574	356	605	60	2798	574	577	584	362	62
28	27	8.24 Premium	J	VS1	62.5	57	355	397	384	247	608	57	2768	573	578	35	62	55	2881	574	571	356	611	60	2798	607	61	572	61	
29	28	8.30 Premium	J	VS2	65.3	58	356	403	407	263	617	61	2769	581	587	348	605	58	2881	604	836	358	615	61	2798	581	581	581	362	61
30	29	8.23 Very Good	J	VS2	60.5	61	357	396	387	24	613	56	2769	582	586	357	62	59	2881	615	847	402	633	55	2798	584	588	588	368	62
31	30	8.23 Very Good	J	VS1	60.9	57	357	396	388	242	625	59	2770	585	581	352	628	58	2882	602	846	364	584	56	2798	587	587	587	362	60
32	31	8.23 Very Good	J	VS1	60	57	402	4	403	241	614	59	2770	583	585	358	618	56	2882	523	236	324	625	55	2798	571	585	585	61	61
33	32	8.23 Very Good	J	VS1	59.8	57	402	404	408	242	628	57	2770	579	581	388	59	58	2882	597	844	351	645	59	2798	628	621	403	62	
34	33	8.23 Very Good	J	VS1	60.7	58	402	397	401	242	614	57	2770	579	588	362	606	58	2882	636	832	36	627	59	2798	61	622	368	61	
35	34	8.23 Very Good	J	VS1	59.8	58	402	401	408	24	607	56	2770	577	578	351	587	57	2882	587	848	349	642	52	2799	581	577	577	371	61
36	35	8.23 Very Good	J	VS1	60.3	58	402	401	408	24	607	56	2770	582	587	358	609	58	2882	587	848	349	642	52	2799	581	577	577	371	61
37	36	8.23 Good	J	VS1	58.2	59	402	406	408	237	615	58	2770	579	575	355	627	64	2883	609	6	379	604	59	2799	582	584	584	366	6
38	37	8.23 Good	J	VS1	64.1	59	402	383	385	248	593	57	2770	582	588	354	606	61	2883	601	838	360	618	58	2799	582	586	586	361	60
39	38	8.23 Very Good	J	IS1	64	58	403	426	431	262	614	57	2770	582	587	348	605	55	2884	632	627	611	608	58	2799	648	648	648	369	60
40	39	8.26 Very Good	J	VS2	60.8	59	403	413	416	252	603	56	2771	581	583	351	630	57	2885	586	6	370	61	60	2799	587	601	608	369	60
41	40	8.23 Very Good	J	VS2	60.3	58	403	408	407	242	625	59	2771	585	581	352	628	58	2885	602	846	364	584	56	2799	587	587	587	362	60
42	41	8.31 Ideal	J	IS2	61.2	58	403	449	45	275	618	57	2771	573	577	358	615	55	2885	538	833	364	609	57	2799	528	53	322	63	
43	42	8.30 Good	J	IS1	60.3	58	403	449	455	276	618	55	2771	587	602	368	603	63	2885	622	841	346	601	61	2799	584	588	588	364	59
44	43	8.30 Good	J	VS2	60.8	58	403	408	402	248	618	56	2771	588	588	358	621	58	2885	602	846	364	584	56	2799	584	588	588	364	59
45	44	8.30 Good	J	IS1	60.4	58	403	419	424	248	61	57	2772	601	602	367	614	58	2886	579	579	358	609	60	2799	585	582	573	371	64
46	45	8.30 Good	J	IS2	61.2	58	403	407	407	248	618	56	2772	587	587	357	614	58	2886	588	846	364	584	56	2799	585	582	573	371	64
47	46	8.30 Premium	Z	IS1	62.4	58	403	424	428	265	617	56	2772	533	537	33	655	55	2875	63	925	411	636	57	2799	641	627	418	61	
48	47	8.30 Very Good	J	IS2	61.8	55	403	435	442	271	602	61	2772	58	588	346	603	61	2875	608	833	401	613	58	2799	584	588	588	364	59
49	48	8.30 Very Good	J	IS2	60.8	58	403	426	438	278	613	56	2772	582	587	348	605	55	2884	632	627	611	608	58	2799	648	648	648	369	60
50	49	8.26 Very Good	J	VS2	63.3	60	404	4	403	244	608	62	2773	587	578	338	622</													

Introduction

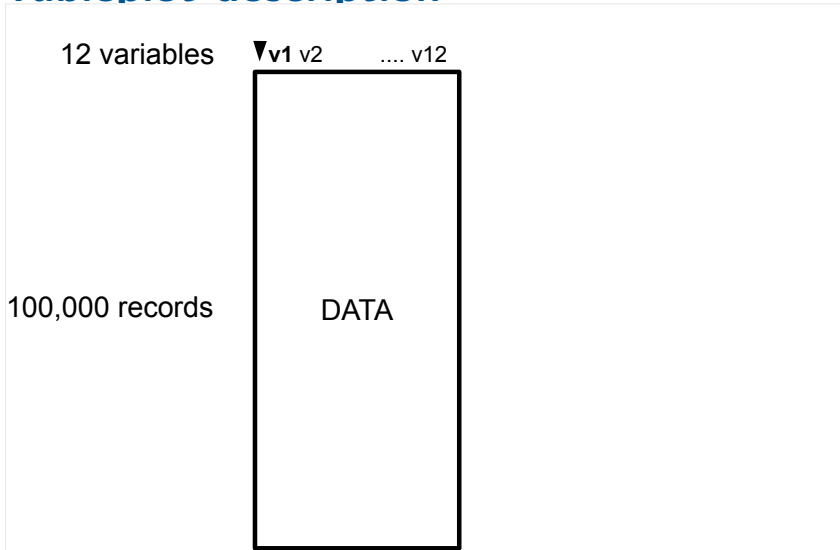
New approach:



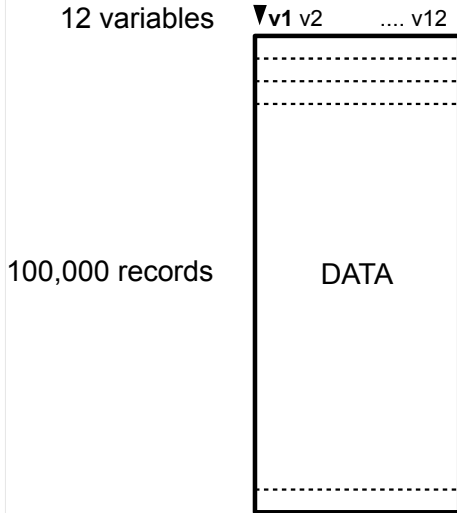
Tableplot description



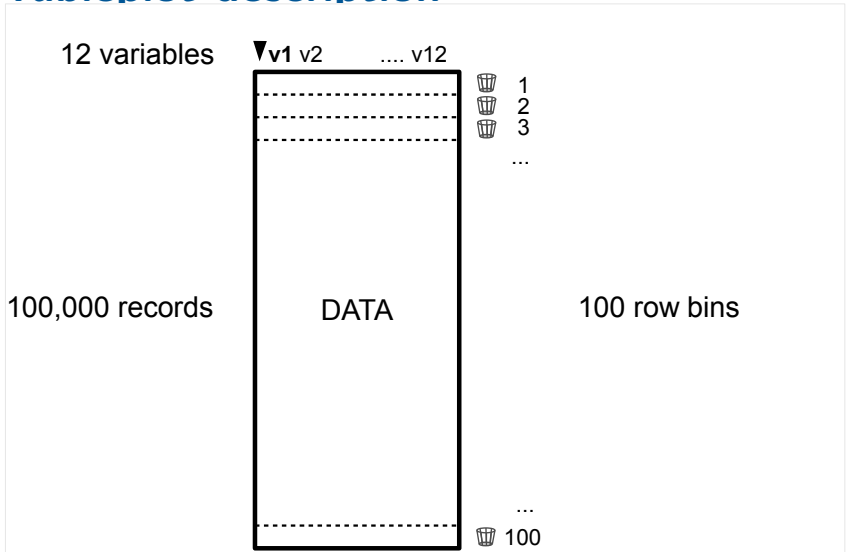
Tableplot description



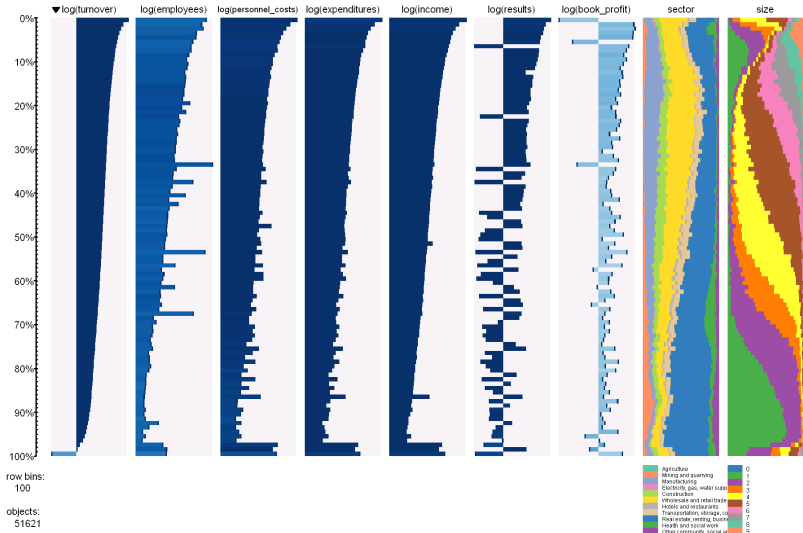
Tableplot description



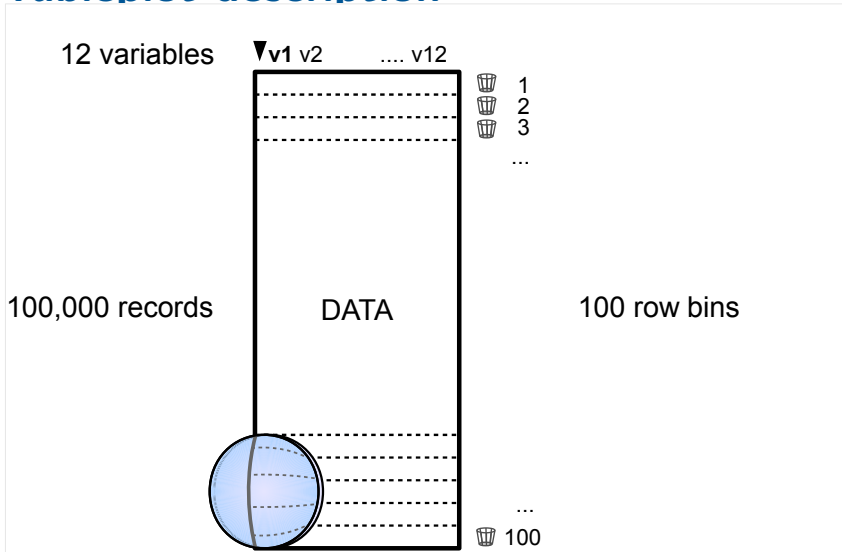
Tableplot description



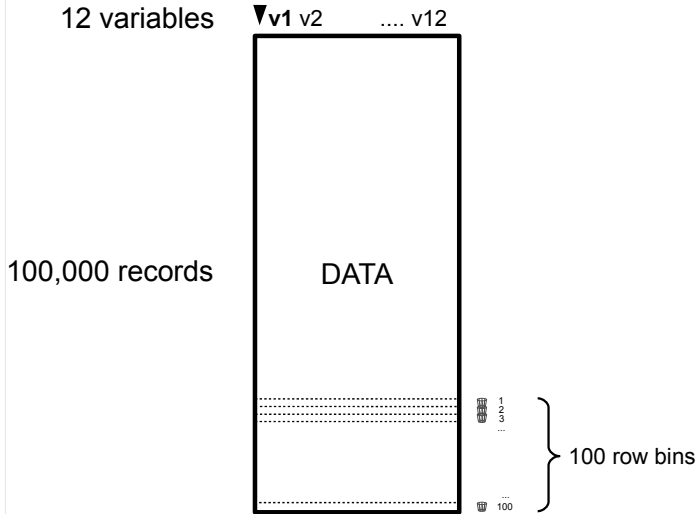
Tableplot description



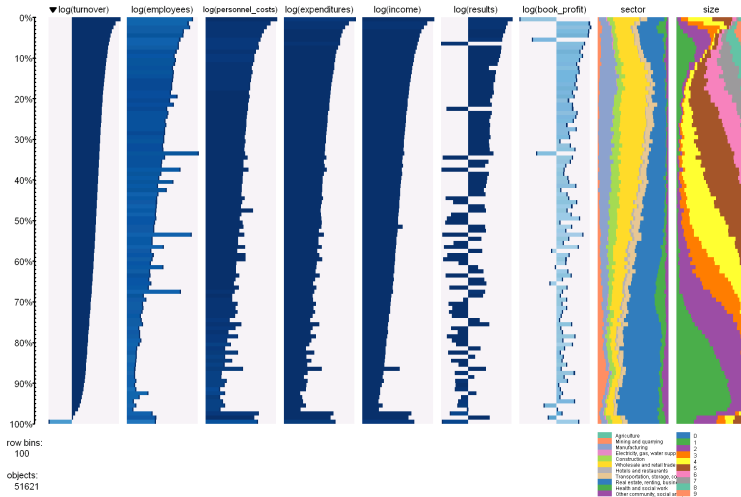
Tableplot description



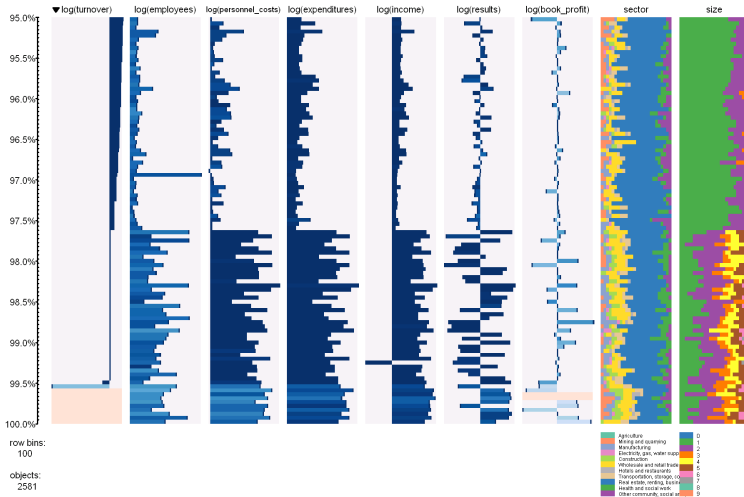
Tableplot description



Tableplot description



Tableplot description



Tableplot description

Quality measures:

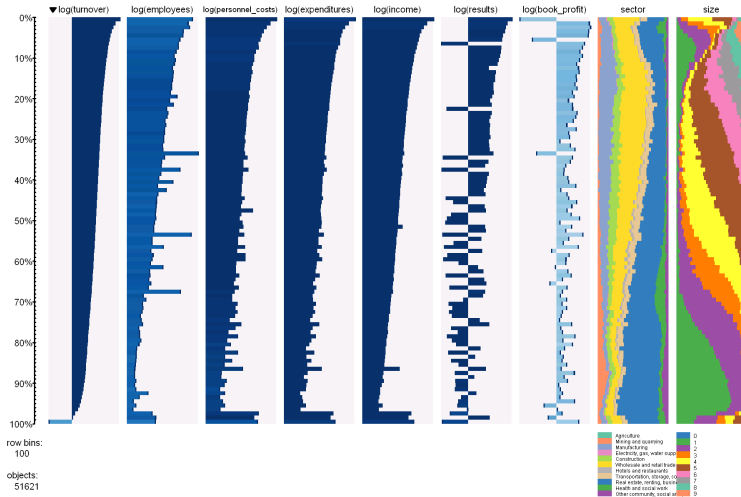
- 1 Smoothness of a data distribution
- 2 Selectivity of missing values
- 3 Distribution of correlated variables

Applications

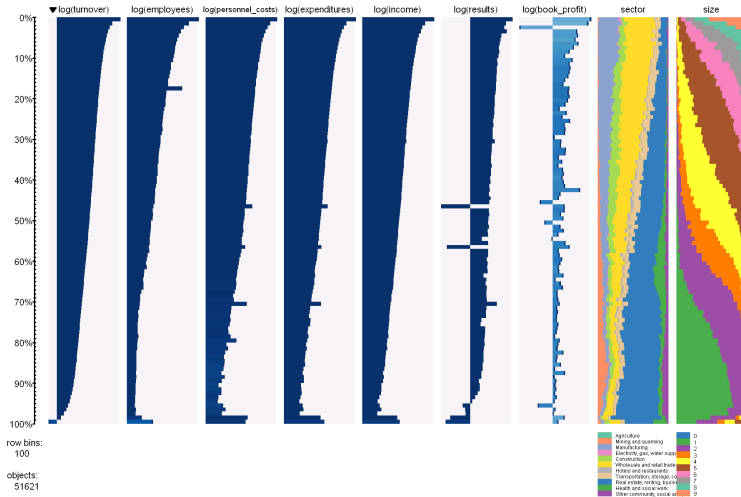
Structural Business Statistics (SBS)

- Large business survey
- Circa 50,000 respondents
- Data editing and analysis process:
 - 1 Unprocessed data
 - 2 Edited data
 - 3 Data prepared for analysis

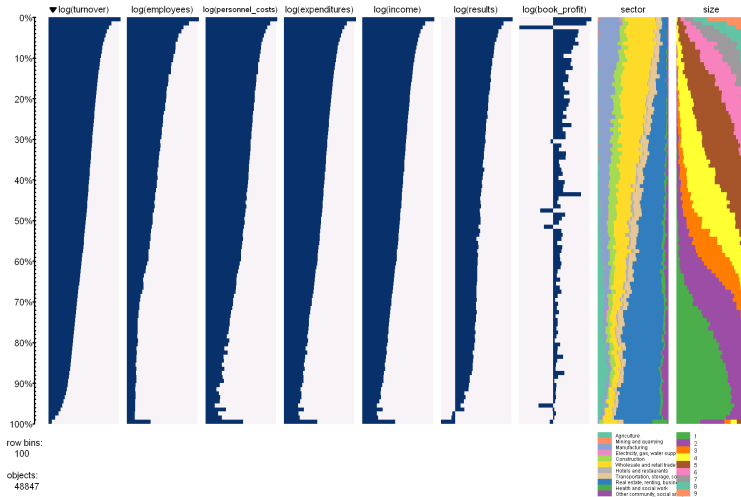
Unprocessed data



Edited data



Data prepared for analysis

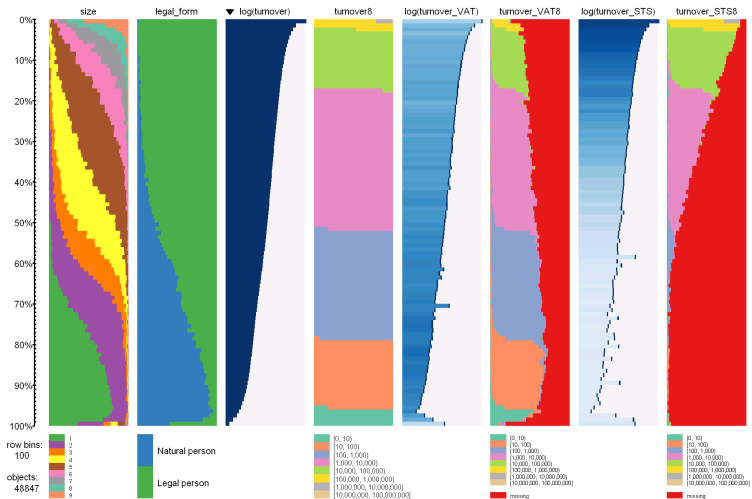


Comparison with other sources

Comparison:

- SBS turnover
- VAT turnover
- STS turnover

Comparison with other sources



Implementation in R

- Package **tabplot**
- Available on CRAN
- Functions

```
tableplot tableplot(myDataFrame,  
                    colNames = myColumnNames,  
                    nBins = 100)
```

```
num2fac num2fac(myNumericVector,  
               method="pretty",  
               n=5)
```

```
tabGUI tableGUI()
```

- Supports very large datasets (up to 2.10^9 records)

Conclusion

- Quality assessment
 - Existing data sources
 - New data sources
- Effective method to support top-down data analysis
- Apply tableplot to other sources
- Further improve tableplots

Conclusion

- Quality assessment
 - Existing data sources
 - New data sources
- Effective method to support top-down data analysis
- Apply tableplot to other sources
- Further improve tableplots

Conclusion

- Quality assessment
 - Existing data sources
 - New data sources
- Effective method to support top-down data analysis
- Apply tableplot to other sources
- Further improve tableplots

Conclusion

- Quality assessment
 - Existing data sources
 - New data sources
- Effective method to support top-down data analysis
- Apply tableplot to other sources
- Further improve tableplots